

Renibacterium salmoninarum:
Genome Sequencing, Finishing
Strategies and Assembly
Validation

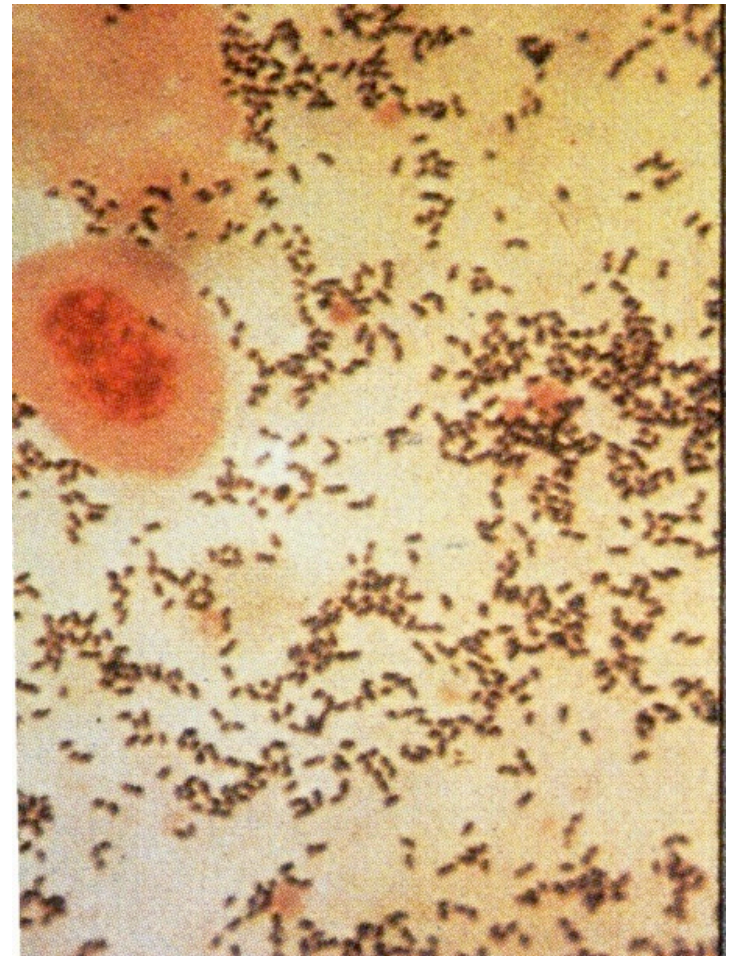
Rajinder Kaul, Ph.D.
University of Washington
Genome Center

Microbial Genomics Era

- First microbial genome sequenced in 1995
Haemophilus influenza RD strain
- Number microbial genomes listed : 827
- Completed: 281
 - archae: 25
 - bacteria: 257
- In Sequencing: 546
 - archae: 25
 - bacteria: 521
- UWGC microbial WGS
 - Completed: 7 organisms
 - Total microbial genome size: 22 Mb
 - In sequencing pipeline: 23 Mb

Renibacterium salmoninarum

- Gram+ve diplococcobacillus, member of the *Micrococcaceae* family
- Grows slowly *in vitro*
 - optimum temperature 15°C
 - doubling time 24 hours
- Infects salmonids,
 - transmitted horizontally and vertically
 - causes bacterial kidney disease (BKD)
- 3.5 Mb genome



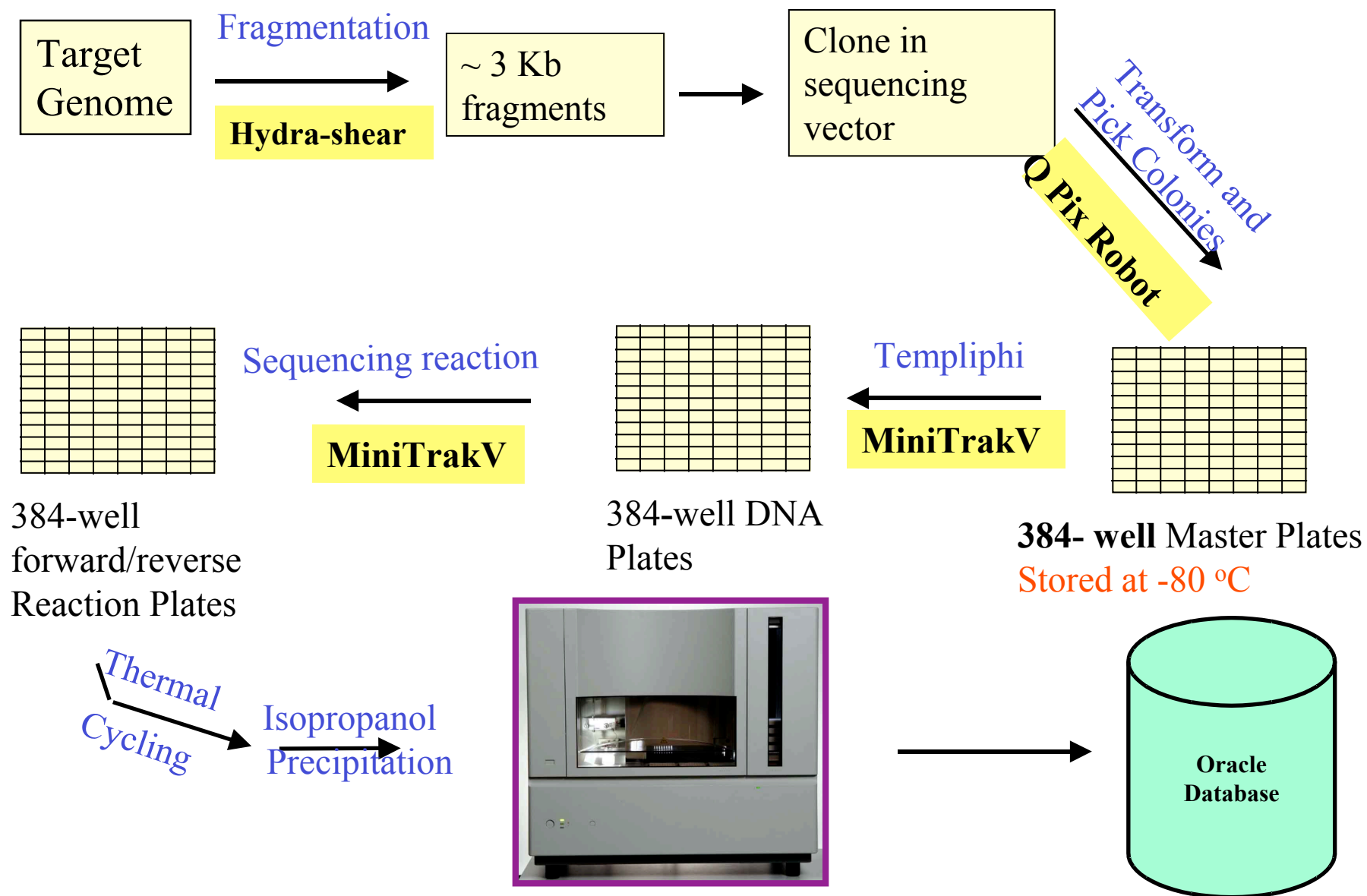
Courtesy: Mark Strom

BKD Incidence

- Predominant salmonid bacterial disease
- 30-85% incidence in hatchery stocks and wild stocks
- Up to 80% mortality in captive stocks
- 15% mortality in captive brood stocks due to BKD since 2000

Courtesy: Mark Strom

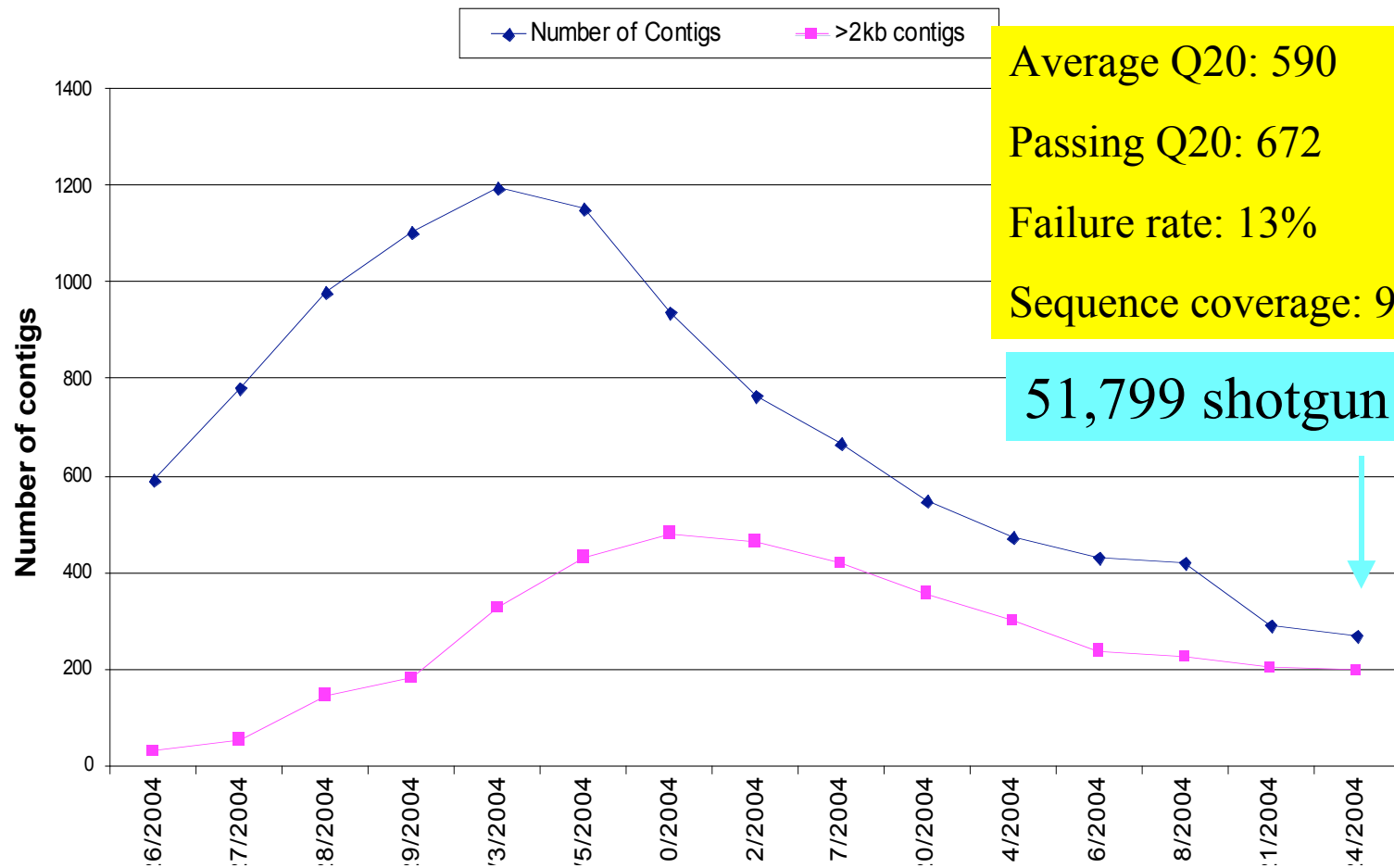
Schematics of Sequencing Pipeline



Some Terminologies

- PHRAP: Software tool for pair-wise alignment of individual reads for assembling in to contiguous chromosome
- PHRED quality: Statistical quality value assignments to trace data
- CONSED: Software tool for viewing genome assemblies and its manual manipulations
- AUTOFINISH: Automated evaluation of sequence quality and contiguity and design experiments to improve quality and attempts to close gaps
- Q20: Logarithmic scale error probability value. Q20 implies probability of 1 in 100 a base could be wrong

Sequencing Progress of Renibacterium salmoninarum



Average Q20: 590

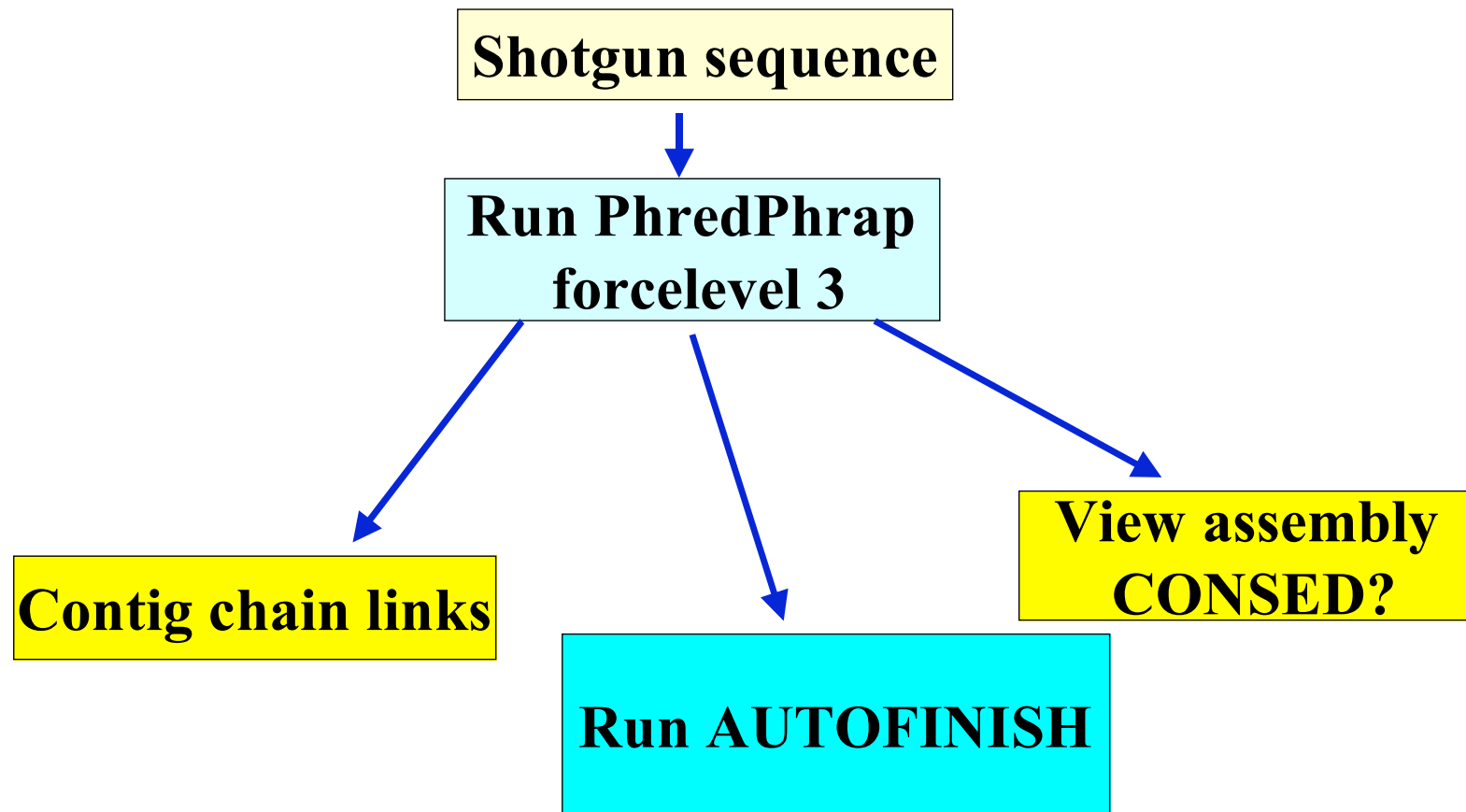
Passing Q20: 672

Failure rate: 13%

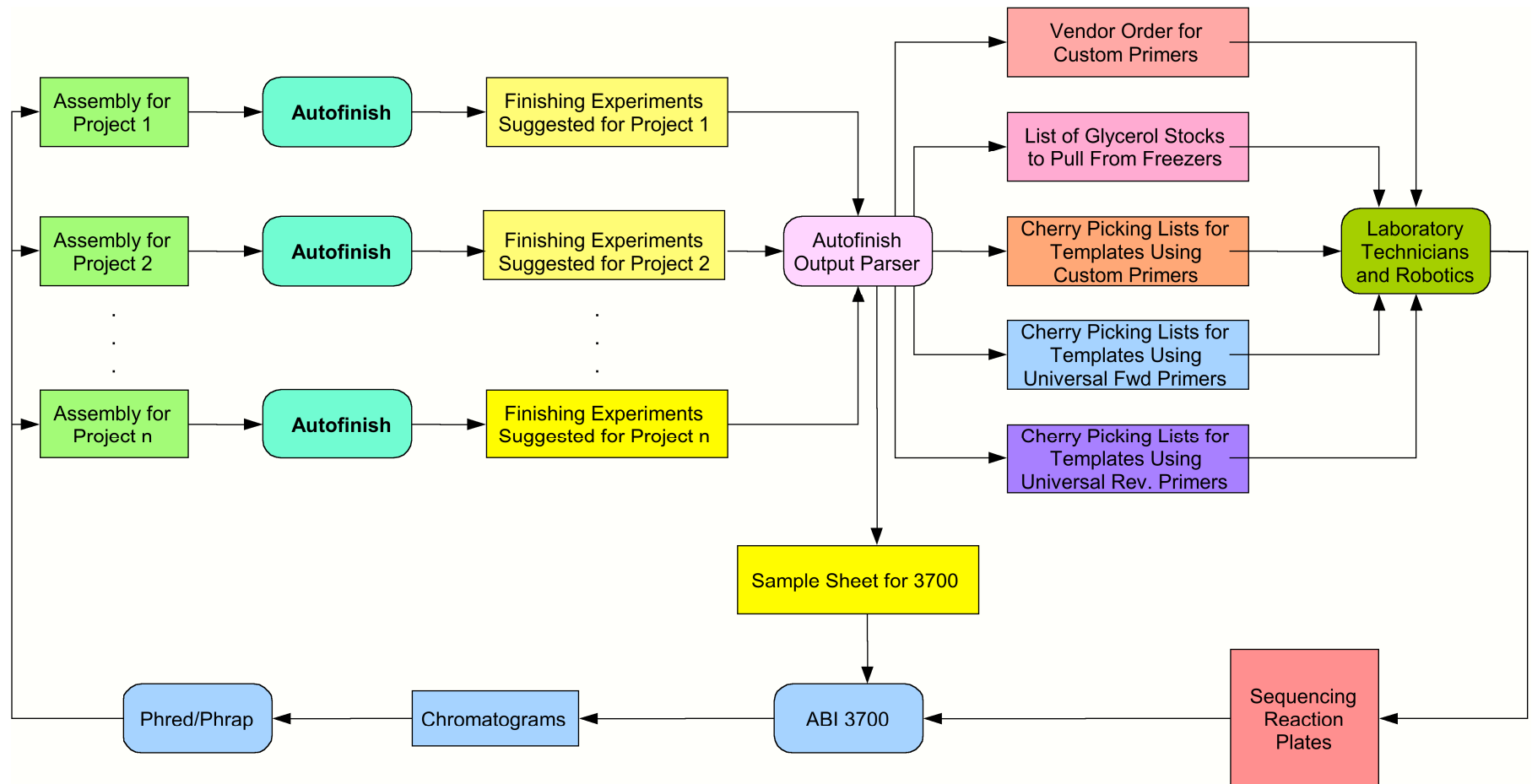
Sequence coverage: 9.7 X

51,799 shotgun reads

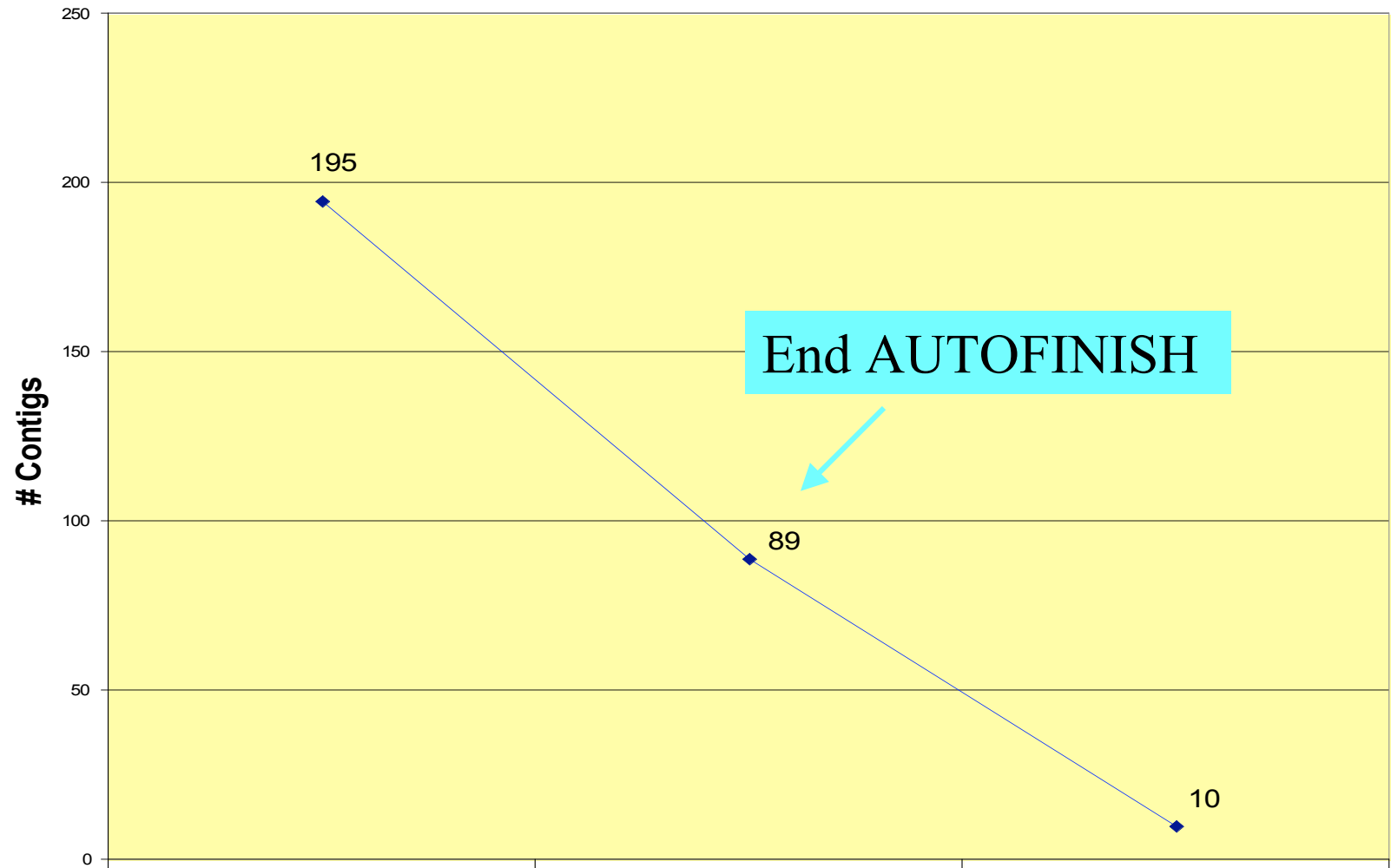
Project Evaluation Following Shotgun Sequencing



Automated Finishing Dataflow



Progress in Finishing



The screenshot displays the Geneious 6.1.10 software interface, showing a sequence alignment of reads against a reference contig (Contig91). The top panel displays the alignment with various reads (e.g., rs33209a2, rs33209a3) and their corresponding sequence positions. The bottom panel shows a detailed view of the alignment for a specific read (rs33209a2_fp55q464.x) with a chromatogram plot below it. The chromatogram shows peaks for each base (A, C, G, T) across the sequence positions. The interface includes a menu bar (File, Navigate, Info, Color, Dim, Misc), a toolbar, and a status bar at the bottom.

gapture - X-Win32

Aligned Reads

File Navigate Info Color Dim Misc Help

rs33209.fasta.screen.ace.34 Contig92 Sone Tags Pos: 15320 cClear

Search for String Compl Cont Compare Cont Find Main Win Err/10kb: 0.18

15300 15310 15320 15330 15340 15350 15360 15370 15380

CONSENSUS GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp08q178.x1u14 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp08q178.x1u14_a GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp08q178.x1u16_b GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp08q178.x1u16_ed GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp08q178.x1u1_e1088 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp20q421.x1 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp20q421.x1r1p341_e1971 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp20q421.x1u14 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp20q421.x1u14_a GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp20q421.x1u16_ed GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp20q421.x1u2_e3597 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp20q421.x2r1p341_e1971 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp20q421.x2u2_e3597 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp21q334.y1 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp20q421.x1r1p341_e1971 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp21q404.x1r2p993_e4669 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp21q404.y1 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp22q114.x1r2p993_e4670 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp22q114.y1 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp22q194.x1r3p1200_e5303 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp22q194.y1 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp22q454.x1r3p1200_e5304 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp22q454.x2r3p1200_e5304 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp22q454.y1 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp29q167.y1 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp29q167.y1u16_b GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp29q167.y1u1_e1965 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp29q167.y2u14 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp29q167.y2u14_a GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp29q269.y1 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp29q269.y1u16_b GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp29q269.y1u1_e1966 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp42q264.x1r2p807_e4067 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp42q264.y1 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp42q264.y1u1_e1093 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp45q471.x1r4p1431_e5997 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp49q277.x1 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp49q277.x1u2_e4061 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp49q277.x2u2_e4061 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp49q435.x1 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp49q435.x1r1p199_e1104 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a2_fp49q435.x2r1p199_e1104 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a3_fp02q132.x1 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a3_fp02q132.x1r2p1036_e4838 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a3_fp02q132.x1r2p704_e3601 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a3_fp02q132.x2u2_e4063 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

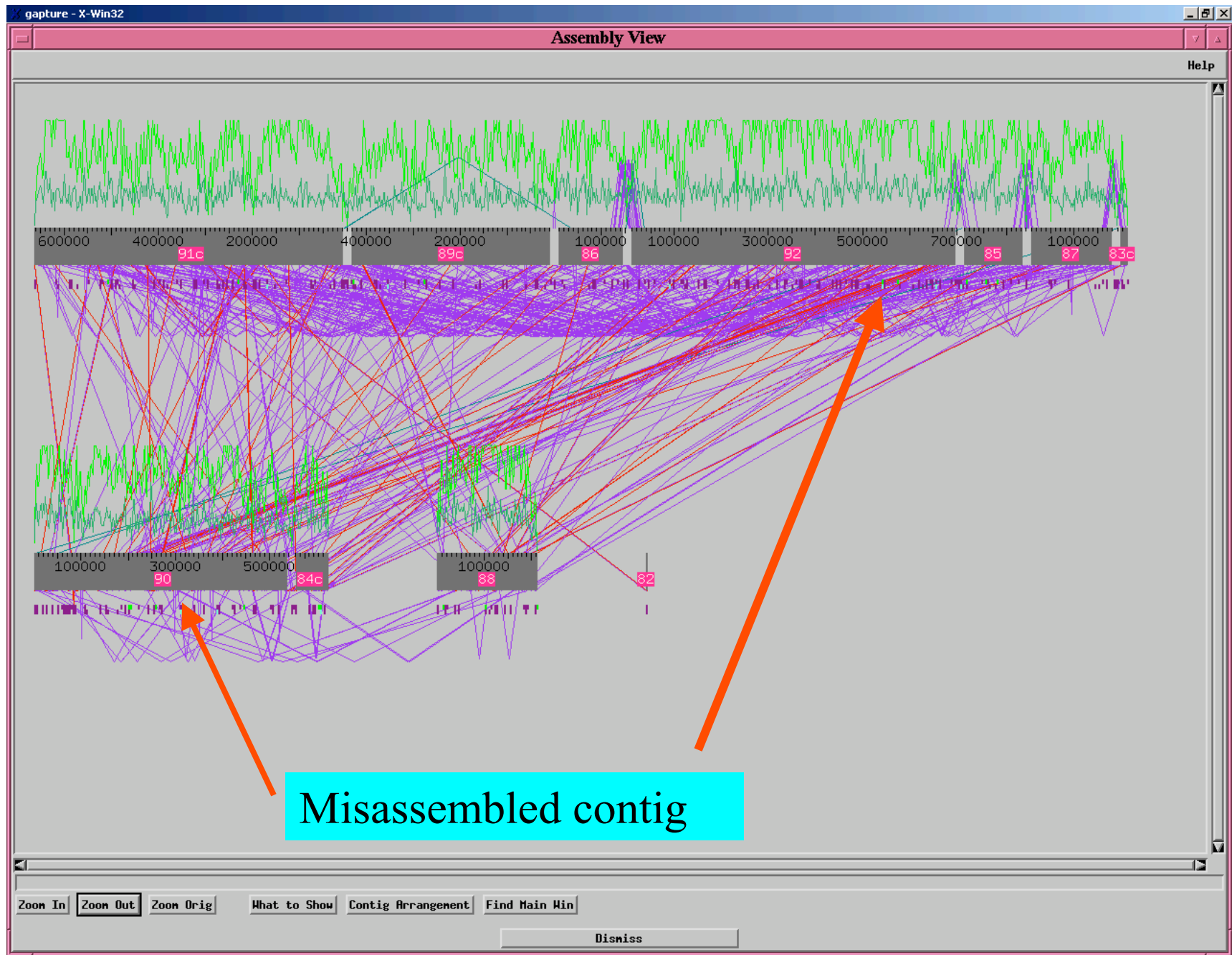
rs33209a3_fp02q132.y1 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a3_fp02q132.y1u2_e4667 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

rs33209a3_fp04q302.x1 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

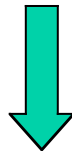
rs33209a3_fp04q302.x1r2p704_e3602 GCG*AGCC*TACA*ACGAGTTTTCGAGCT**ACTCT*CTGATCCTTCGTTACGGGA*TTAAACGAAATCG*CC*CTC*GATATTGGTTGGAGCCGGAGAC

<<< << < Prev Next > >> >>> cursor dismiss

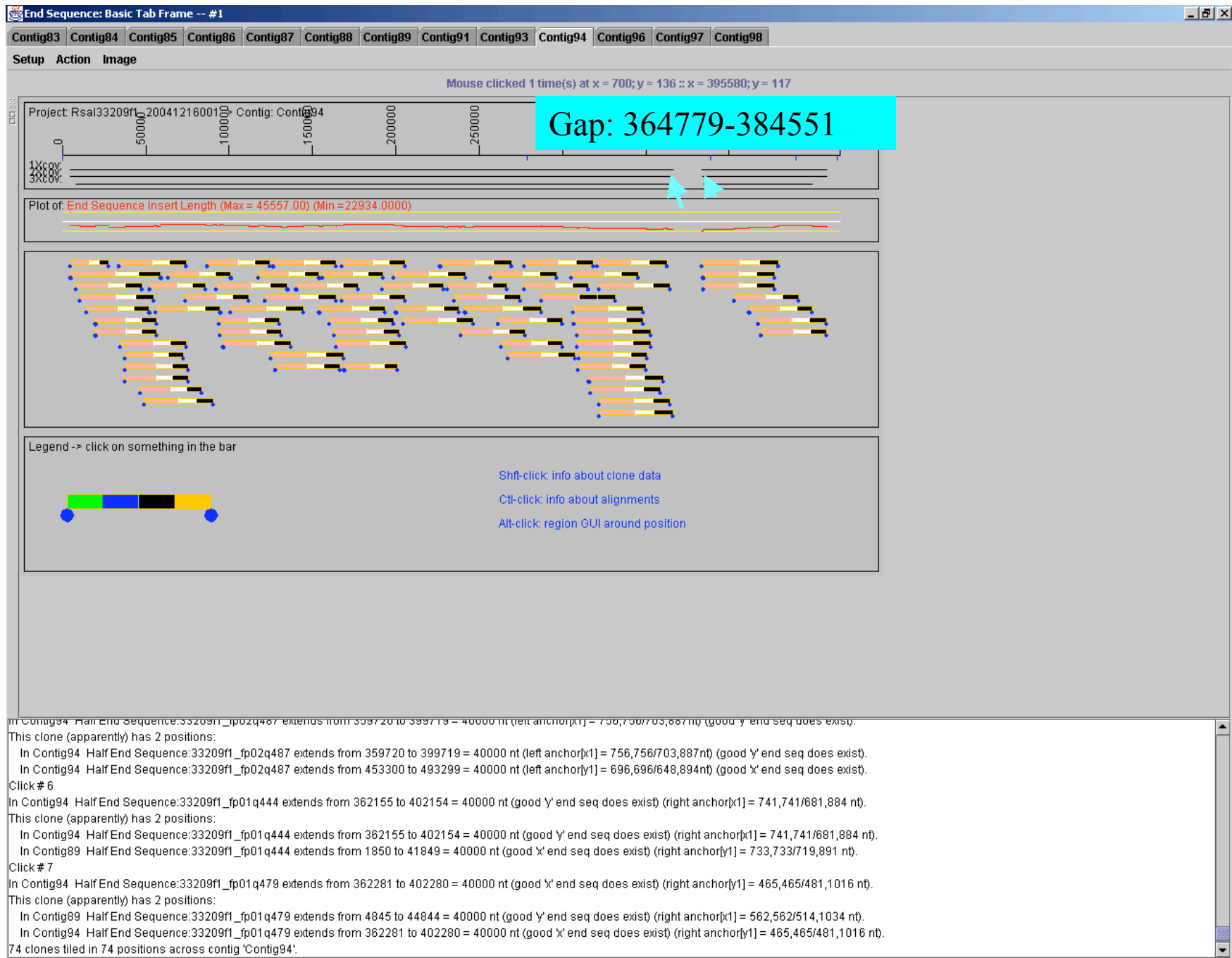


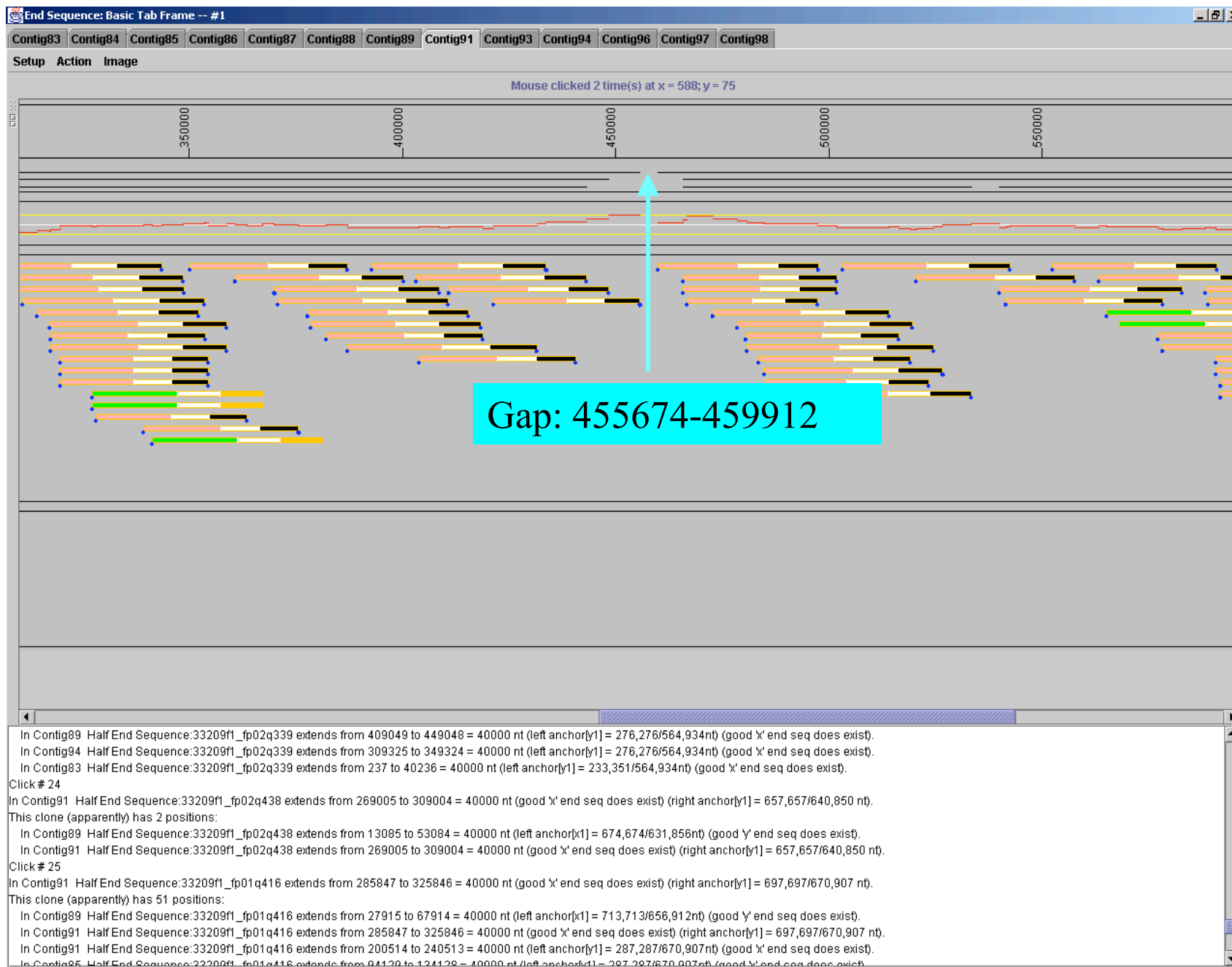
Assembly Validation Genomic Variation

Make fosmid libraries from genome

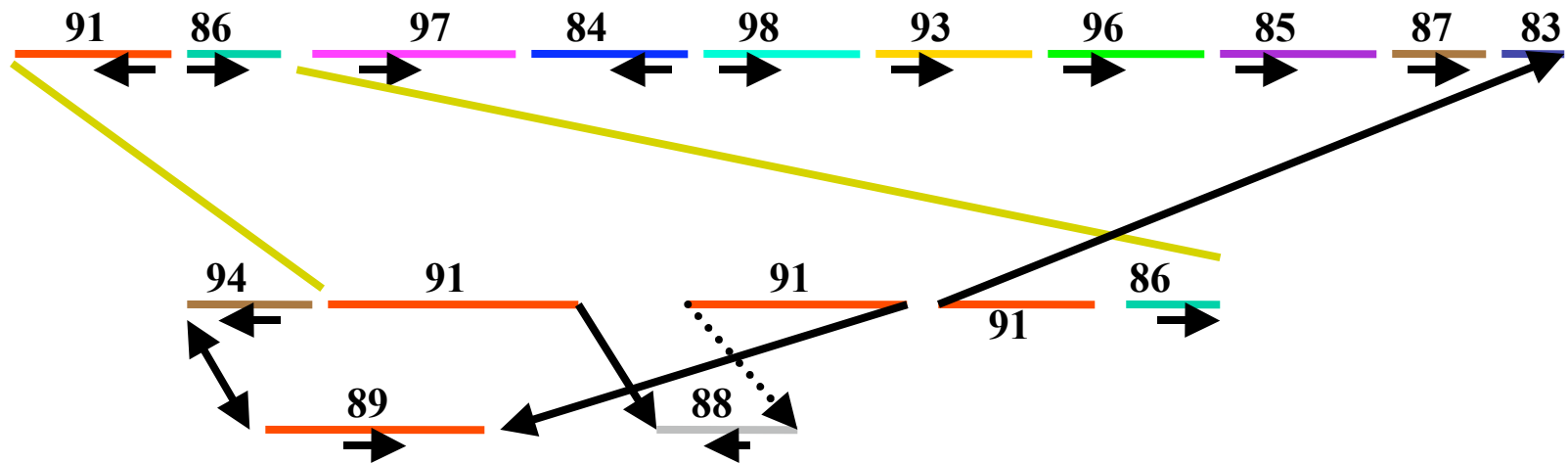


End Sequence and Fingerprint data used
to validate genome assembly





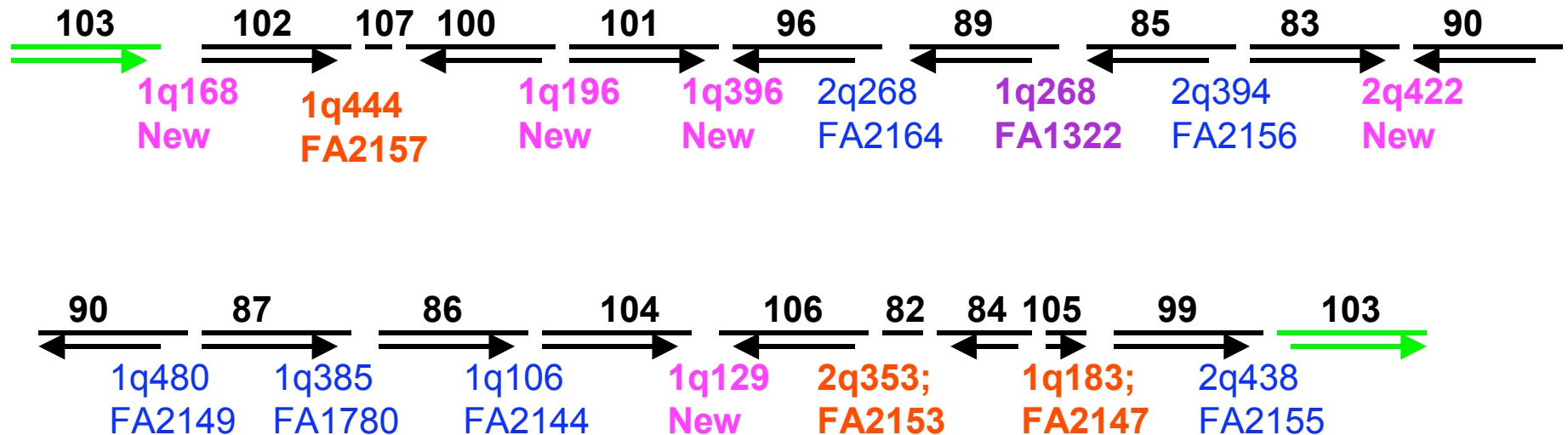
Contig Order and Orientation



- 11 original contigs
- Fosmid end sequence tiled pointing misassembled regions
- Identify fosmids to fix misassemblies

Manual Tearing and Order-Orientation

Renibacterium salmoninarum Ace.40; 7/13/2005



- Tore contigs to create 20 from 11 original contigs
- Individual fosmids were sequenced
- Backbones from fosmid assemblies used to fix misassemblies
- Low quality regions improved by additional sequencing

Renibacterium salmoninarum

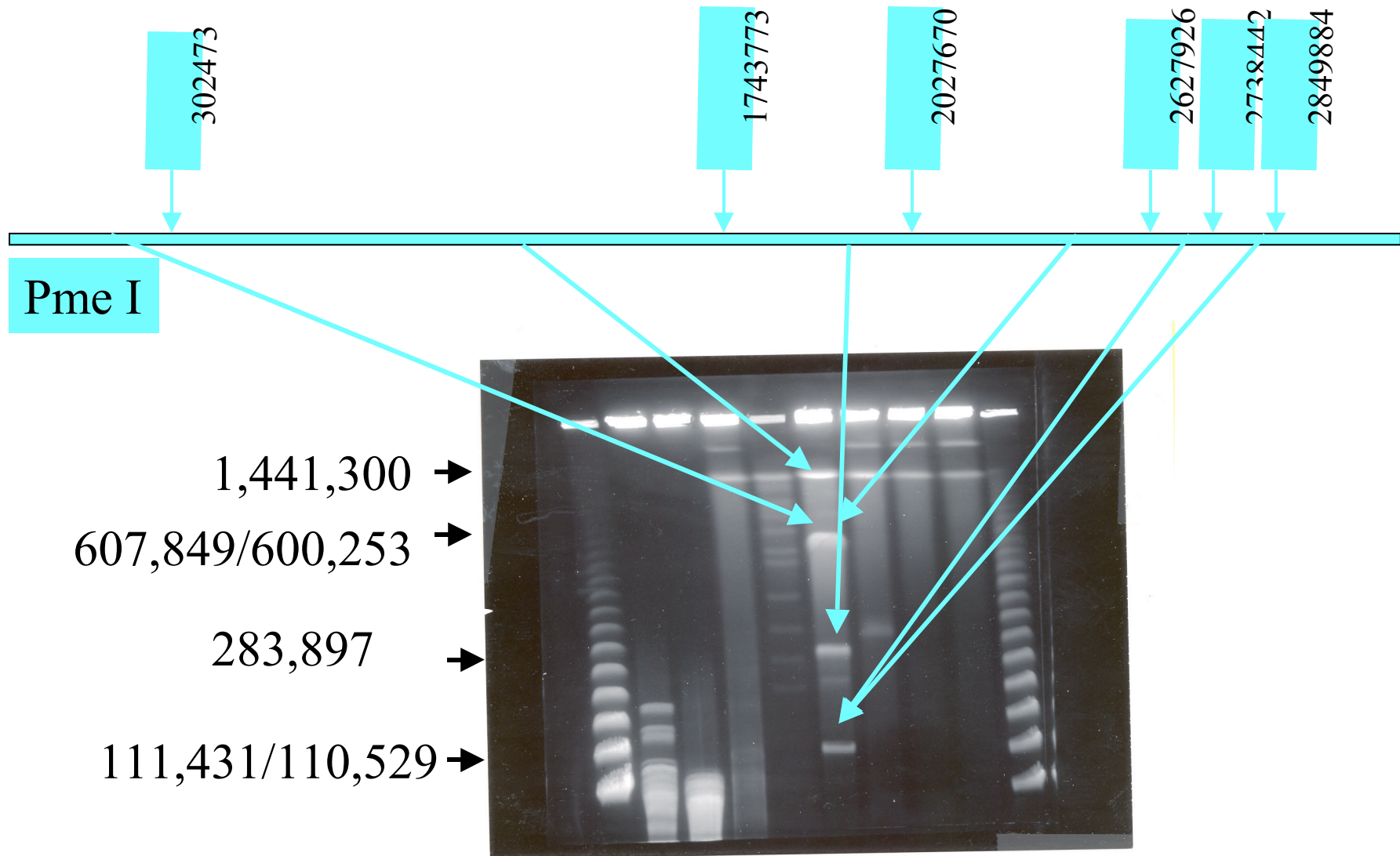
Genome

- Single Circular chromosome with 3,155,258 bases
- 52692 reads in final assembly
- 51799 shotgun reads
- 8968 finishing reads
- 29 fosmid clones sequenced independently to fix misassemblies (~22,272 reads)
- 83,932 total reads

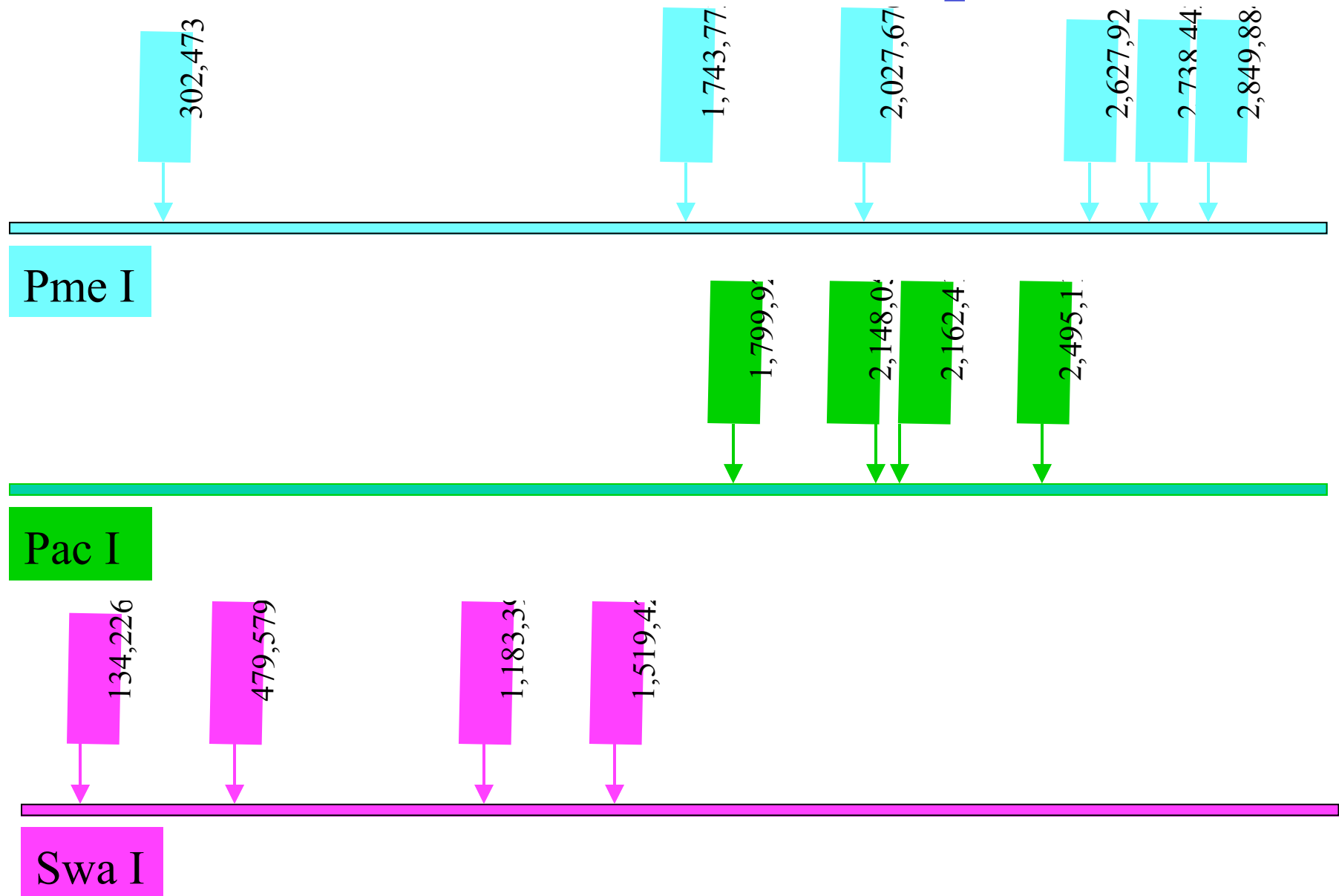
Genome Assembly Validation

- Gross level validation
 - Compare experimental and virtual sequence derived PFGE fragment distribution data derived from rare cutting restriction enzymes
- Finer scale assembly validation (1 Kb resolution)
 - Tile fosmid end-sequences and corresponding fingerprint data in Hind III, Pst I, and Bgl II domains against finished genome assembly
 - All sequence derived virtual fingerprint fragments should be accounted for by experimentally derived data, and vice versa
 - Requires at least 2X fingerprint fragment coverage in 2 domains at any given point in assembly

Pme I: PFGE & Sequence Fragments



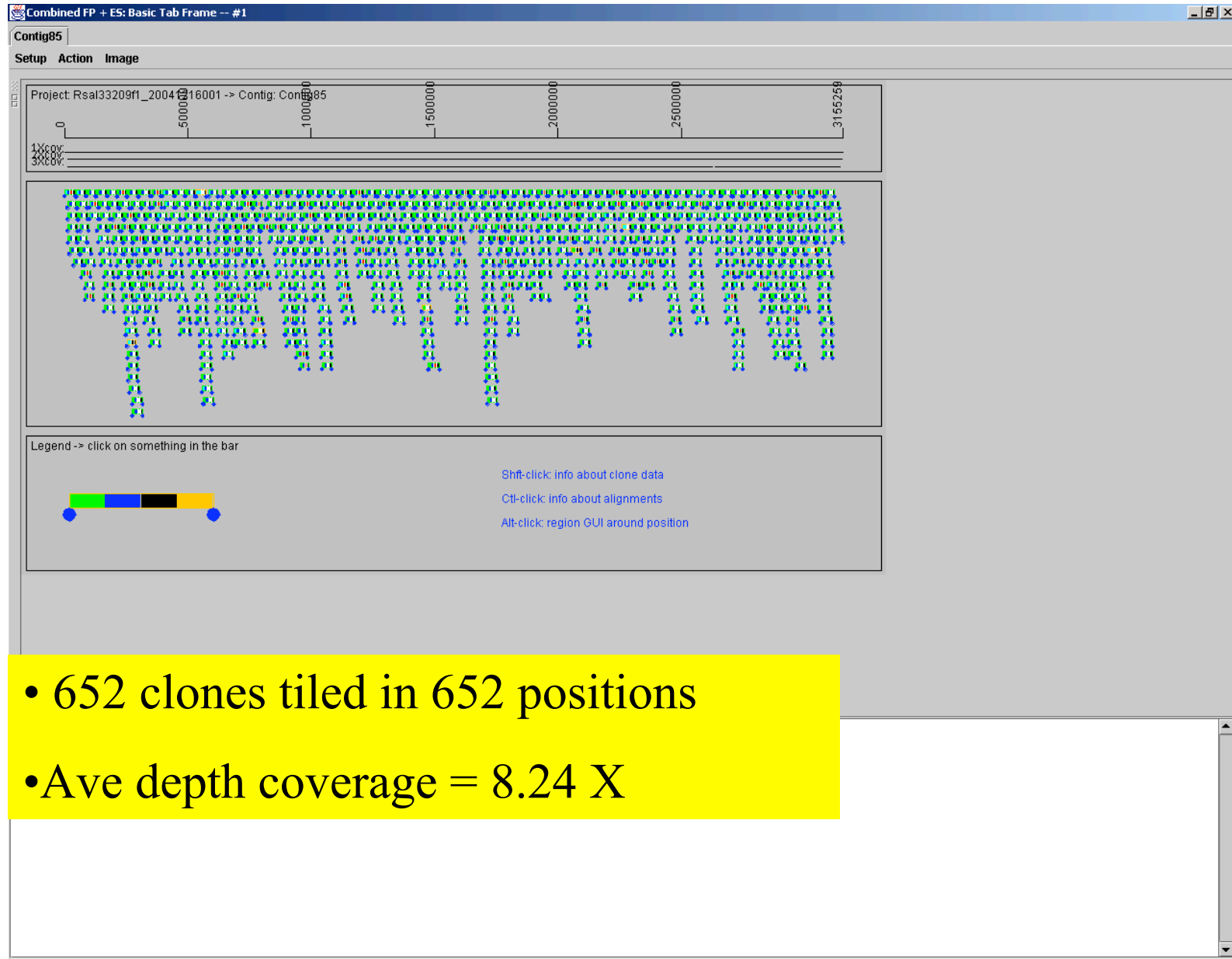
Pulse Field Gel Electrophoresis



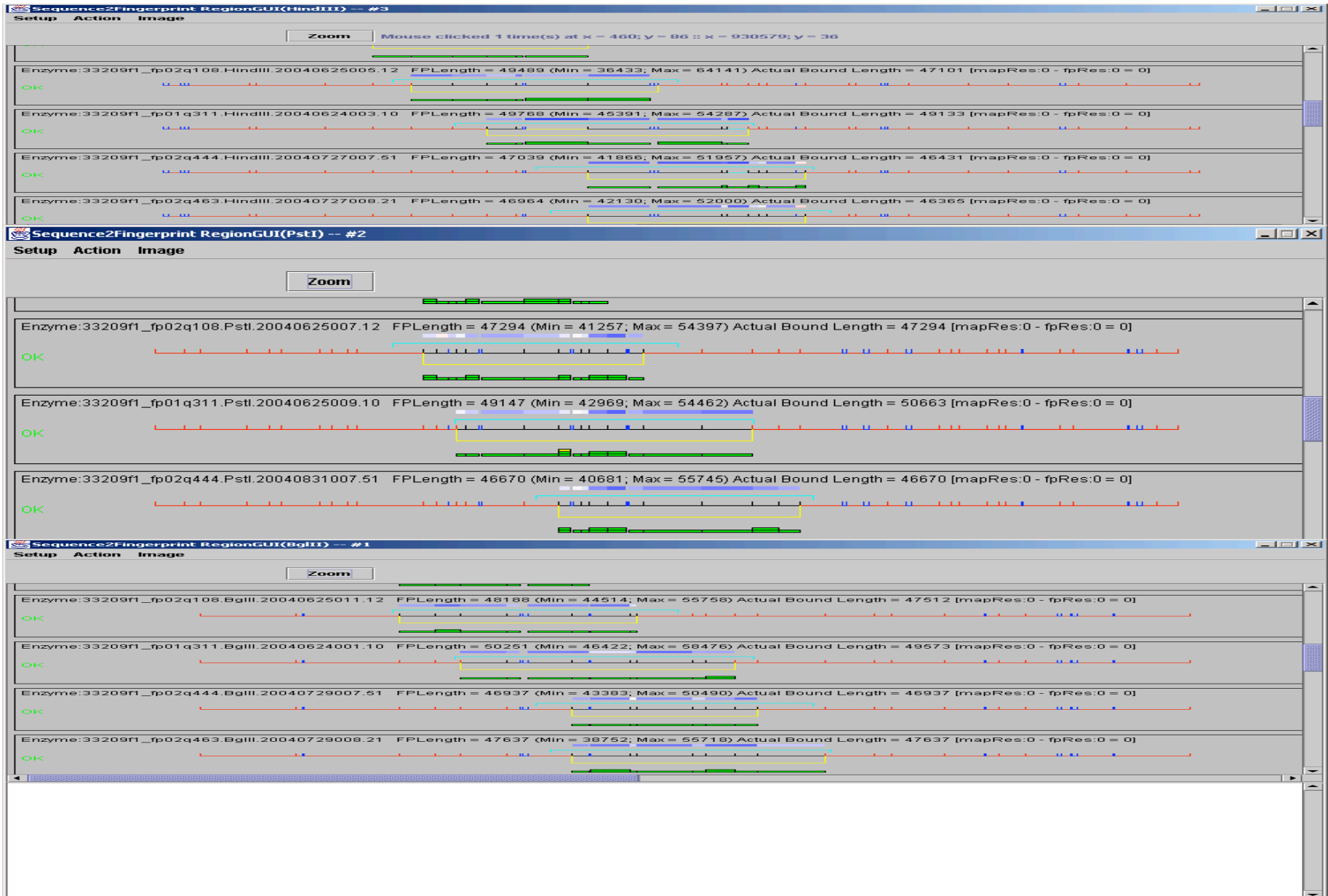
Fosmid End-sequences and FP data

- **Clones = 768 (9.7 X clone coverage)**
 1. clones with X end seq = 728
 2. clones with Y end seq = 725
 3. clones with NO end seq = 20
- **Clones with 3 enzymes present = 730**
 1. clones with BglII = 755
 2. clones with PstI = 754
 3. clones with HindIII = 747
 4. clones with (exactly) 0 enzymes present = 4
 5. clones with (exactly) 1 enzymes present = 2
 6. clones with (exactly) 2 enzymes present = 32

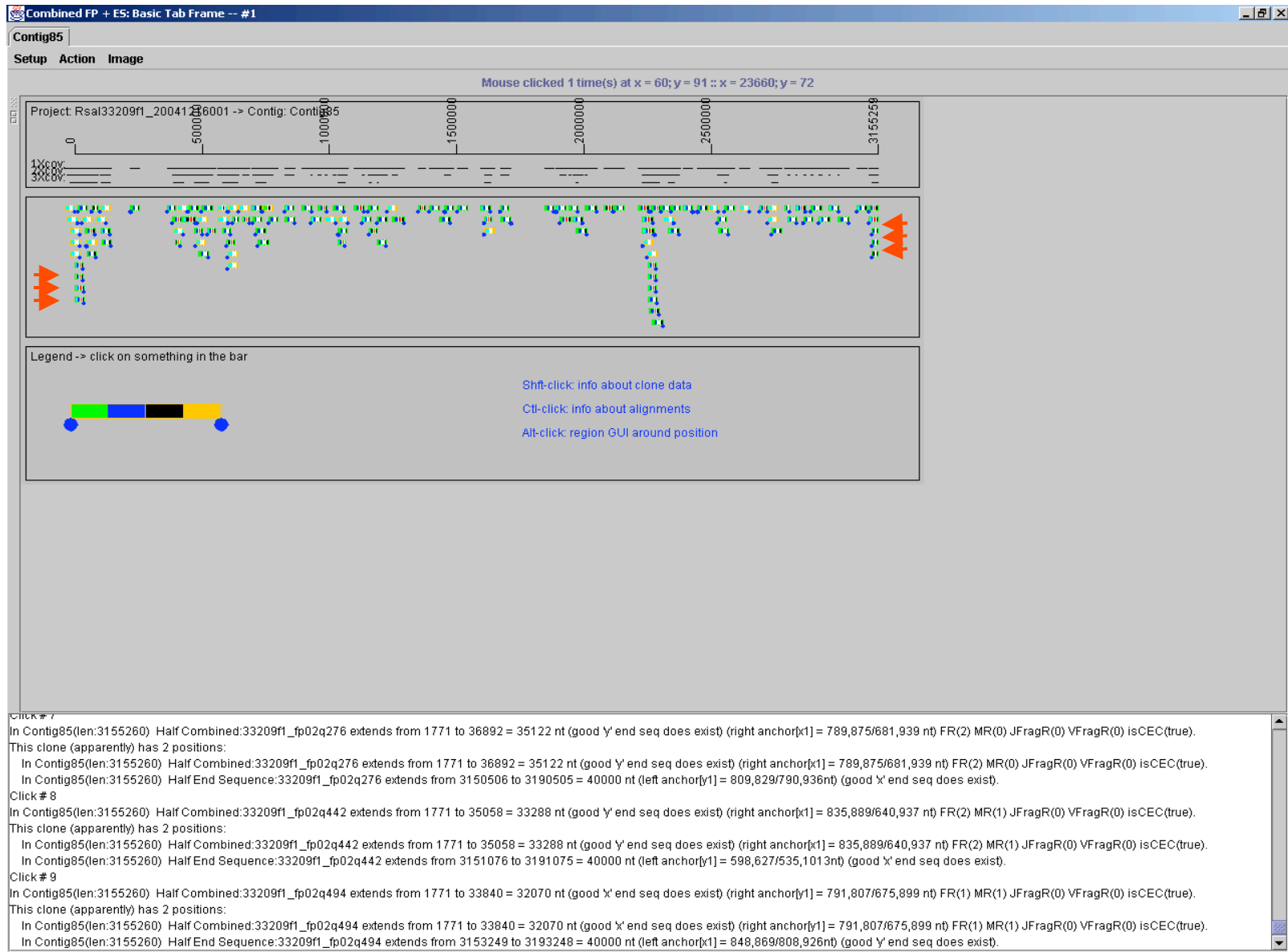
Best Full Position Tiling



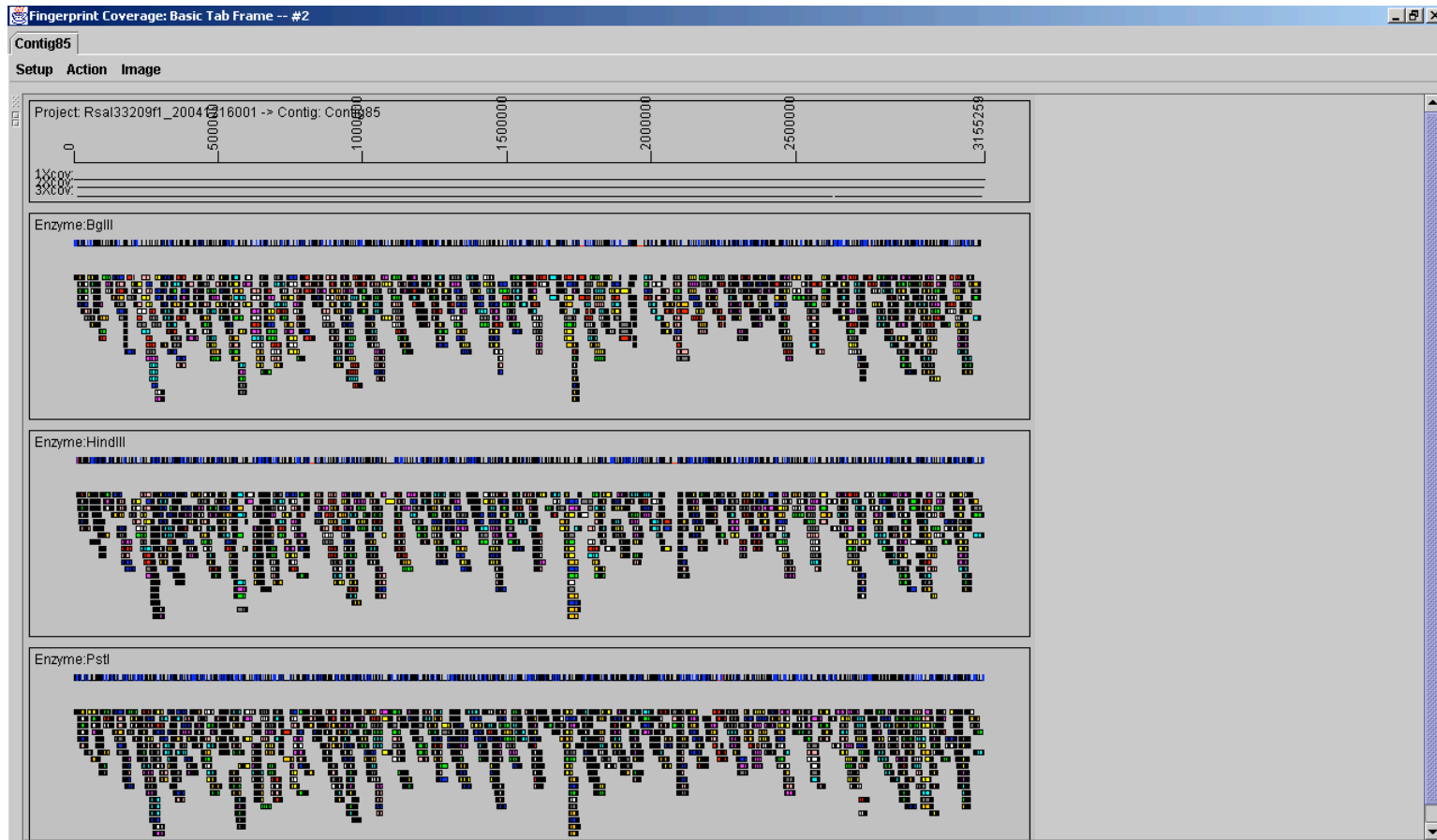
Restriction Fragment Coverage in Multiple Domains Across Clones



Clones Validating Circular Chromosome



Fingerprint Coverage

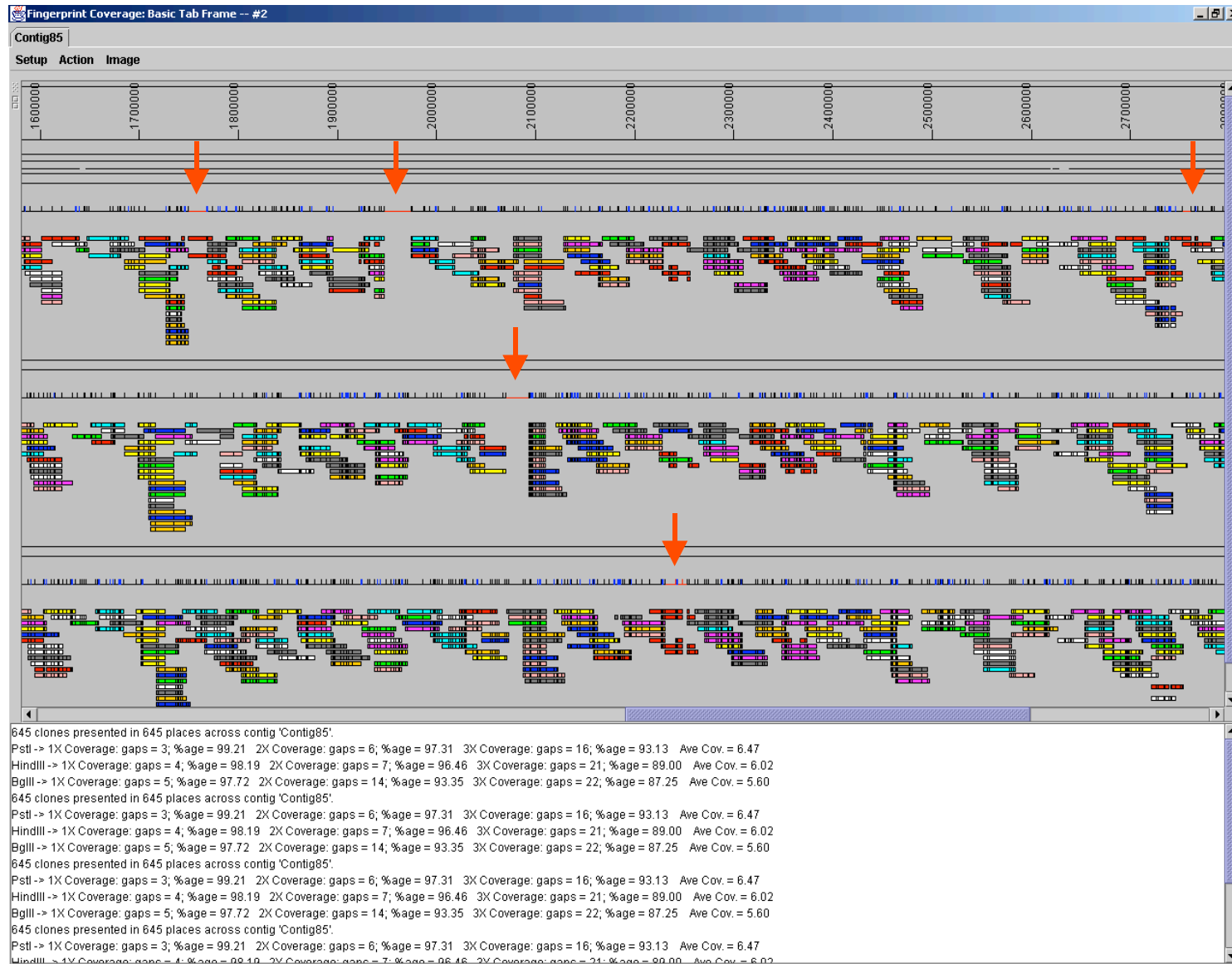


PstI -> 1X Coverage: gaps = 3; %age = 99.21 Ave Cov. = 6.47

HindIII -> 1X Coverage: gaps = 4; %age = 98.19 Ave Cov. = 6.02

BglII -> 1X Coverage: gaps = 5; %age = 97.72 Ave Cov. = 5.60

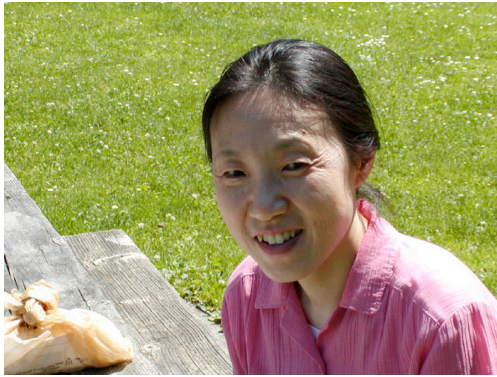
Non-contiguous Gaps in FP Coverage



Genome Sequencing Highlights

- Single Circular chromosome with **3,155,258** bases
- Two insertional sequence repeat families
 - IS994: 70 copies
 - ISRsa2: 11 copies
- GC content: 56.27%
- Number of ORFs: 3,667
- Sequence with coding potential: 91%
- Average ORF size: 862 bases

Acknowledgements



UWGC Production Staff